## NGS: Technology and Current Applications

### SWGDAM
### Quantico, VA

**Dr. Peter M. Vallone**
**Leader, Applied Genetics Group**
**NIST**
**January 7, 2014**

---

## Disclaimer

- Forensic DNA research conducted at NIST is supported by an interagency agreement between the National Institute of Justice and the NIST Law Enforcement Standards Office.

- Points of view in this document are those of the authors and do not necessarily represent the official position or policies of the U.S. Department of Commerce. Certain commercial equipment, instruments, and materials are identified in order to specify experimental procedures as completely as possible.

- In no case does such identification imply a recommendation or endorsement by NIST, nor does it imply that any of the materials, instruments, or equipment identified are necessarily the best available for the purpose.

---

## Disclaimer

- The NIST talks today are intended for educational purposes

- Technology is moving at a fast pace

- If your favorite platform, application, library prep, software, etc. is not mentioned
  **Please bring it up!!!**

## Outline

- Introduction
- Non-forensic applications
- Generalized workflow
- Platforms and throughput
- Sequencing chemistries
- Wrap up / thoughts

## What's in a name???

Massively parallel sequencing

**NGS**

*Second-generation sequencing*

Next-generation sequencing

Whole-genome sequencing

Third-generation sequencing

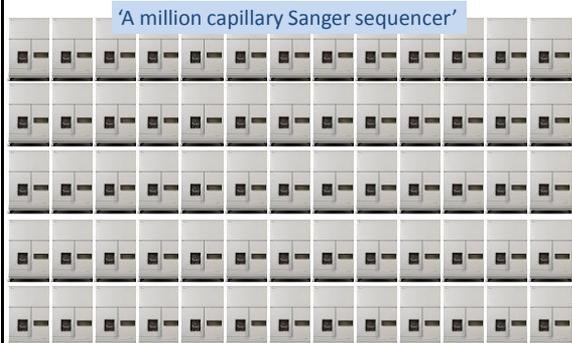**HIGH-THROUGHPUT SEQUENCING**

Next-generation genomics

## Parallel Sequencing

'A million capillary Sanger sequencer'

## Parallel Sequencing

'A million capillary Sanger sequencer'

- Clonal vs population amplification
- Shorter reads (Range 75 to 400)
- Errors are more 'detectable'
- High coverage 100 – 1000 - 10,000x
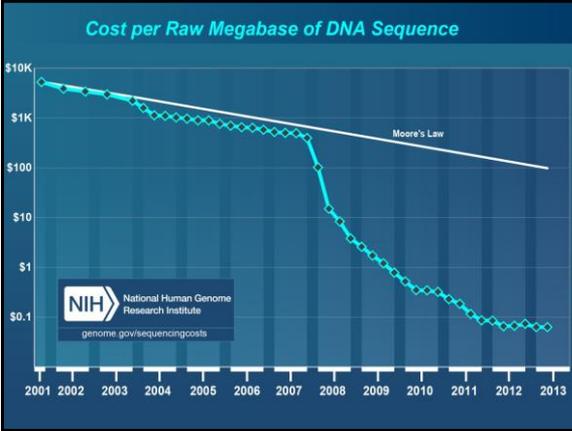- **Rely more on informatics to assemble millions of short reads**

Size 10 font
8.5 x 11 paper
5,580 bases per page
3.234 Gb = 579,570 pages

_As of Nov 2013_

Web of Science
Articles and Reviews
"next generation sequencing"

_As of Nov 2013_

Web of Science
Articles and Reviews
"next generation sequencing" **and**

Cancer
Clinical
Microbial
Virus
Metagenomics
Prenatal
Transplant
Forensic
STR
Forensics

## Non-forensic applications

- Clinical
- Inherited disease
- Reproductive health
- Cancer – gene fusion
- Rare variants
- Pre-implantation (genetic screening)
- Transplant medicine (HLA)
- Microbiomics/Metagenomics
- Gene expression | RNA seq
- Public health
- Ancient DNA
- NIPT (non-invasive prenatal testing)

## Non-forensic applications

- Scanning Nature Reviews Genetics

## Generalized NGS Workflow

≈500-1000 ng of genomic DNA — *Library preparation* — Hours to days

Genomic DNA → Fragment to ≈200 bp → Ligate PCR adapters → PCR → Sequencing

One template per bead/droplet/spot

Illumina

PGM "Ionogram"

## Generalized NGS Workflow

≈1-5 ng of genomic DNA

*Library preparation*

Hours to days

Genomic DNA → Fragment to ≈200 bp → Ligate PCR adapters → PCR → Sequencing

alternative

One template per bead/droplet/spot

PCR amplicons ← Clean up and quant PCR products

Illumina

Target specific genes or regions
CODIS STRs, SNPs
≈500 ng of **PCR product**

PGM "Ionogram"

---

## Whole Genome versus Targeted

- Whole genome
  - Genomic DNA sheared and sequenced
    - 500-1000 ng of DNA template

- Targeted
  - PCR amplified or hybridization captured regions of the genome are sheared and sequenced
    - Start with 1-5 ng of DNA -> amplify/enrich to 500-1000 ng

---

## Generalized NGS Workflow

Minutes - Hours

Sequencing → Sequence reads → Assemble → Evaluate Variants

Mb to Gb of data

To reference sequence or *de novo*

coverage

FASTQ format

SAM/BAM format (aligned to a reference)

## FASTQ Format

- FASTQ - normally uses four lines per sequence.

- Line 1 begins with a '@' character and is followed by a sequence identifier and an *optional* description
- Line 2 is the raw sequence letters.
- Line 3 begins with a '+' character and is *optionally* followed by the same sequence identifier (and any description) again.
- Line 4 encodes the quality values for the sequence in Line 2, and must contain the same number of symbols as letters in the sequence.

- `(1)@SEQ_ID`
  `(2)GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCC`
  `(3)+`
  `(4)!''*((((***+))%%%++)(%%%%).1***-+*''))**`

http://maq.sourceforge.net/fastq.shtml

## Aligning Sequencing Reads

- One common algorithm is BWA
  - Burrows-Wheeler Aligner
  - Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler Transform. Bioinformatics, 25:1754-60.
  - Li H. and Durbin R. (2010) Fast and accurate long-read alignment with Burrows-Wheeler Transform. Bioinformatics, Epub.

**Considerations when choosing an alignment software**
- Speed
- Memory
- Accuracy
- Variant calling (SNPs, InDels)

## SAM/BAM Format

- SAM Sequence Alignment Map Format
- Simple, tab-delimited text file
- BAM (optional compressed binary encoding)

| | Type | | | |
|---|---|---|---|---|
| | Field | Regular expression | Tag | Range | Description |
| Alignment | QNAME | [^ \t\n\r]+ | | | Query pair NAME if paired; or Query NAME if unpaired [2] |
| | FLAG | [0-9]+ | | [0,2^16-1] | bitwise FLAG (Section 2.2.2) |
| | RNAME | [^ \t\s\r#=]+ | | | Reference sequence NAME [3] |
| | POS | [0-9]+ | | [0,2^29-1] | 1-based leftmost POSition/coordinate of the clipped sequence |
| | MAPQ | [0-9]+ | | [0,2^8-1] | MAPping Quality (phred-scaled posterior probability that the mapping position of this read is incorrect) [4] |
| | CIGAR | ([0-9]+[MIDNSHP])+|\* | | | extended CIGAR string |
| | MRNM | [^ \t\n\r#]+ | | | Mate Reference sequence NaMe; "=" if the same as <RNAME> [3] |
| | MPOS | [0-9]+ | | [0,2^29-1] | 1-based leftmost Mate POSition of the clipped sequence |
| | ISIZE | -?[0-9]+ | | [-2^29,2^29] | inferred Insert SIZE [5] |
| | SEQ | [acgtnACGTN.=]+|\* | | | query SEQuence; "=" for a match to the reference; n/N/. for ambiguity; cases are not maintained [6,7] |
| | QUAL | [!-~]+|\* | | [0,93] | query QUALity; ASCII-33 gives the Phred base quality [6,7] |
| | TAG | [A-Z][A-Z0-9] | | | TAG |
| | VTYPE | [AifZH] | | | Value TYPE |
| | VALUE | [^\t\n\r]+ | | | match <VTYPE> (space allowed) |

http://samtools.sourceforge.net/SAMv1.pdf
http://chagall.med.cornell.edu/NGScourse/SAM.pdf

## CLC bio Sequence Viewer



## CLC bio Sequence Viewer

Position 750 in the rCRS
A -> G Transition



## Variant Call Table



Reference Position
Type
Length
Reference
Allele
Zygosity
Count
Coverage
Frequency
Forward-reverse balance
Average quality

## Variant Call Table

Rows: 55   Variants                                                                    Filter

**Quality-based Variant Detection (Tue Nov 19 16:28:44 EST 2013)**
**Version:** CLC Genomics Workbench 6.5.1
**User:** ngs
**Parameters:**
```
        Neighborhood radius = 5
        Maximum gap and mismatch count = 2
        Minimum neighborhood quality = 15
        Minimum central quality = 20
        Ignore non-specific matches = Yes
        Ignore broken pairs = Yes
        Minimum coverage = 10
        Minimum variant frequency (%) = 1.0
        Maximum expected alleles = 1
        Advanced = No
        Require presence in both forward and reverse reads = No
        Ignore variants in non-specific regions = Yes
        Filter 454/Ion homopolymer indels = No
        Create track = Yes
        Create annotated table = Yes
        Genetic code = 1 Standard
```
**Comments:** Edit
Found 86 variants (including reference alleles)
**Originates from:**
    9947a_S2_L001_R1_001_2 (paired) (Reads) (history)

## Platforms

- Illumina
  - **MiSeq**
  - HiSeq 2000/2500
  - GAIIx
- Life Technologies
  - SOLiD (5500 series)
  - **Ion Torrent PGM**
  - Ion Torrent Proton
- Pacific Biosciences
  - PACBIO RS II
- 454 Roche
  - GS jr
  - GS FLX+

October 15, 2013 – Roche shutting down 454 sequencing business Will be phased out mid-2016

## On the horizon…

- Qiagen GeneReader
  - Sequencing by synthesis approach
  - Should be available in 2014
  - QiaCube NGS (for automated library preparation)
  - Qiagen has also purchased CLC bio
- Oxford Nanopore
  - Ratcheting strand of DNA through a protein manifold
  - Bases are detected by a difference in current

## Moving Targets
6 months from now these parameters will have changed

- Newer instruments
- Costs decreasing
- Throughput increasing
- Read lengths increasing
- Chemistries improving
- Library preparations – simpler/automated
- Computers faster – data storage cheaper
- Platforms leaving the market (e.g. Roche 454)
- Platforms entering the market (e.g. Qiagen GeneReader)

---

## Low Throughput versus High Throughput

|  | Illumina MiSeq | Ion Torrent PGM | PacBio RS | Illumina GAIIx | Illumina HiSeq 2000 |
|---|---|---|---|---|---|
|  | **Benchtop** | | | **High Throughput** | |
| Instrument Cost | $128 K | $80 K | $695 K | $256 K | $654 K |
| Sequence yield per run | 1.5-2 Gb | 100-200 Mb 316 chip | 100 Mb | 30 Gb | 600 Gb |
| Cost/Gb | $502 | $1000 | $2000 | $148 | $41 |
| Run time | 27 hours | 2 hours | 2 hours | 10 days | 11 days |
| Observed raw error rate | 0.80 % | 1.71 % | 12.26 % | 0.76 % | 0.26 % |
| Read length | 150 (300) | 200 (400) | 1500 | 150 | 150 |
| Input DNA | 50-1000 ng | 100-1000 ng | 1 ug | 50-1000 ng | 50-1000 ng |

Adapted from: Quail et al. BMC Genomics 2012, 13:341
http://www.biomedcentral.com/1471-2164/13/341

---

## Balancing the Equation
What question are you trying to answer?

- What instrument and/or strategy is right for my application?

- Markers

- Coverage

- Samples

- Cost (per sample and unit of information)

Other relevant questions:
- Input amounts?
- Desired level of accuracy?
- Integrity of DNA?
- Mixtures present?

## Balancing the Equation
### What question are you trying to answer?

Platform 'X' provides 2 Gb of sequence per run

- Markers
  - 25 STRs and 1000 SNPs
    - 1 STR = 500 bp
    - 1 SNP = 50 bp
- Coverage
  - 600x
- Samples
  - 48

$$=[(25*500)+(1000*50)]*600*48 = 1.8 \text{ Gb}$$

## Balancing the Equation
### What question are you trying to answer?

**MiSeq Output Calculations**

| | MiSeq with:<br>- Upgraded hardware, or from September 2012 and later<br>- MCS v2.3 or later<br>- MiSeq Reagent Kit v3 | MiSeq with:<br>- Upgraded hardware, or from September 2012 and later<br>- MCS v2.0 or later<br>- MiSeq Reagent Kit v2 |
|---|---|---|
| Reads/flow cell | 25,000,000 | 16,000,000 |
| Genome or region size (in bases) | 16,569 | 16,569 |
| Coverage | 10000 | 10000 |
| Total number of cycles (e.g. 300 for 2x150) | 300 | 300 |
| Total output required (in bases) | 165,690,000 | 165,690,000 |
| Output/flow cell (bases/flow cell) | 7,500,000,000 | 4,800,000,000 |
| Number of flow cells | 0.02 | 0.03 |
| Number of samples per flow cell | 45.27 | 28.97 |

The numbers in this spreadsheet are reasonable expectations assuming flow cells are clustered at the proper density. Output may vary based on sample quality, cluster density and other experimental factors. Use these calculations as estimates for planning your runs.

For more information about calculating coverage estimates, see the Coverage Calculation Tech Note.

http://support.illumina.com/downloads/sequencing_coverage_calculator.ilmn

## Multiplexing Samples - Barcoding

- A sample can be tagged with a unique sequence (during library preparation)
- The tagged samples could then be sequenced together and separated in the analysis stage

| Bases in barcode index | Unique Sequence Possibilities |
|---|---|
| N | 4 |
| NN | 16 |
| NNN | 64 |
| NNNN | 256 |
| NNNNN | 1024 |
| NNNNNN | 4096 |

Trade off volume of sequence information for more samples per run

## Life Tech - Ion Torrent - PGM

- Ion Torrent launched in Feb. 2010
- Ion Torrent sequencing employs an analogous technique as pyrosequencing:
  - Emulsion PCR for single copy reactors
  - Non-labeled nucleotide triphosphates are flowed over a bead on a semiconductor surface
- Hydrogen Ion detection
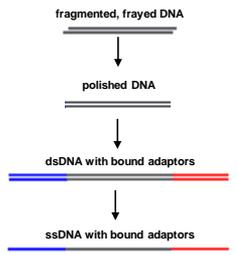  - pH change is detected
  - **No optics**

---

Ion Torrent - PGM

## Constructing a Library

fragmented, frayed DNA

↓

polished DNA

↓

dsDNA with bound adaptors

↓

ssDNA with bound adaptors

- Frayed ends are enzymatically "polished"
- Adaptor oligos ligated onto fragments
- Denatured to ssDNA
- Emulsified with primer-coated beads
- Hybridization of template to bead

*Margulies et al. (2005) Supplementary Materials*

---

Ion Torrent - PGM

## Emulsion PCR & Enrichment

Beads coupled to various library templates

Emulsification & PCR

Oil

- Beads and templates emulsified
- Primer-coated beads bind template
- PCR amplifies template
- Enrich for beads containing PCR products
  - magnetic capture
- Adaptable to automation (Ion Chef)

*Dressman et al. (2003)*

Ion torrent PGM chip

'314' chip

dNTP

Sensing Layer
Sensor Plate

Bulk | Drain | Source | To column receiver
Silicon Substrate

- Chip flooded with one nucleotide after another
- H⁺ released when a complementary base is added to template
- Charge from the ion causes detectable pH change
- Sequencer calls the base

Ion Torrent - PGM

Life Technologies

Sequence: ...AATCTTCTGAATTTCTGCAA....

---

# Illumina MiSeq

- MiSeq launched in Jan. 2011
- The MiSeq uses a sequencing by synthesis approach:
  - Nextera enzymatically fragments and tags DNA
  - Limited cycle PCR
  - Flow cell hybridization
  - Bridge PCR - clusters

- Fluorescent light detection
  - Each base has a unique color
  - Sequence each end of the molecule

---

Illumina - MiSeq

# Nextera Sample Prep/Library Creation

Figure 2: Nextera Sample Preparation Biochemistry

Transposomes

Genomic DNA

~ 300 bp

Tagmentation

~ 300 bp

http://www.illumina.com/documents/%5Cproducts%5Cdatasheets%5Cdatasheet_nextera_dna_sample_prep.pdf

## Topics for further thought

- Additional genetic markers
  - SNPs (ancestry, phenotypic traits, lineage)
  - Insertion Deletion (InDels)
- Data interpretation and review – level of retention
- STR nomenclature
  - Back compatibility with existing databases
  - Future searching methods
- Ethical considerations with coding region markers
- Validation of NGS systems/methods
  - Use of existing standards (SRMs)

## NIST SRM Support

- Further characterization of SRM 2391c, 2392, and 2392-I
- In depth sequencing of mitochondrial genomes and core STR alleles
  - Sanger
  - NGS (PGM and MiSeq)
  - Posters presented at the 25th annual ISFG meeting

"Additional Sequence Characterization of NIST SRM 2391c: PCR-Based DNA Profiling Standard"
http://www.cstl.nist.gov/strbase/pub_pres/Hill-ISFG2013-SRM2391c.pdf
"Characterization of NIST Standard Reference Materials by Next Generation Sequencing"
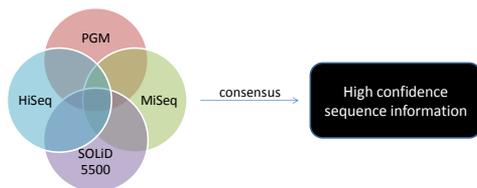http://www.cstl.nist.gov/strbase/pub_pres/KieslerISFG2013poster.pdf

## Multiple NGS Platforms

- Use of multiple platforms to obtain a consensus sequence for the SRMs
  - Identify and reduce the false positives and negatives
  - Identify and control for bias in a specific chemistry and/or informatics pipeline

## Mitochondrial SRMs
### False Positives and False Negatives
Using platform specific informatics pipeline

| | | PGM 1 | PGM 2 | PGM 3 | HiSeq | MiSeq | 5500 |
|---|---|---|---|---|---|---|---|
| 9947A | FP | 1 | 5 | 3 | 21 | 9 | 11 |
| | FN | 3 | 4 | 3 | 3 | 3 | 3 |
| CHR | FP | 2 | 6 | 10 | 21 | 9 | 10 |
| | FN | 3 | 5 | 4 | 3 | 3 | 4 |
| HL-60 | FP | 1 | 8 | 8 | 20 | 9 | 8 |
| | FN | 1 | 2 | 1 | 1 | 1 | 1 |
| Avg Coverage | | 280 | 6,500 | 9,000 | 49,000 | 41,000 | 29,000 |

Calls made to the rCRS
On average 99.94 % agreement with Sanger sequencing

## Heteroplasmy at Position 1,393
### SRM 2392 Component B (9947A)

| Nucleotide Position | rCRS Reference Sequence | SRM 2392 Component B Sanger Call | EdgeBio PGM | NIST PGM run 1 | NIST PGM run 2 | EdgeBio Illumina MiSeq | Beckman Genomics Illumina HiSeq | NIST SOLiD |
|---|---|---|---|---|---|---|---|---|
| 93 | A | G | G | G | G | G | G | G |
| 195 | T | C | C | C | C | C | C | C |
| 214 | A | G | G | G | G | G | G | G |
| 263 | A | G | G | G | G | G | G | G |
| 309.1 | : | C | | | | | | |
| 309.2 | : | C | | | | | | |
| 315.1 | : | C | | | | | | |
| 750 | A | G | G | G | G | G | G | G |
| 1393 | G | G | G/A | G/A | G/A | G/A | G/A | G/A |
| 1438 | A | G | G | G | G | G | G | G |
| 4135 | T | C | C | C | C | C | C | C |
| 4769 | A | G | G | G | G | G | G | G |
| 7645 | T | C | C | C | C | C | C | C |
| 7861 | T | C | C | C | C | C | C | C |
| 8448 | T | C | C | C | C | C | C | C |
| 8860 | A | G | G | G | G | G | G | G |
| 9315 | T | C | C | C | C | C | C | C |
| 13572 | T | C | C | C | C | C | C | C |
| 13759 | G | A | A | | A | A | A | A |
| 15326 | A | G | G | G | G | G | G | G |
| 16311 | T | C | C | C | C | C | C | C |
| 16519 | T | C | C | C | C | C | C | C |

## Heteroplasmy at 1,393?



1,393 G

F873
R2194
F895
F1095
R1769
F1234

*Sequencing primer position*

- 6x coverage by Sanger
- 3/6 of reads indicate low-level heteroplasmy
  – Red circles
- Not reproducible in all reads
  – Not always detected by Sanger sequencing

## Heteroplasmy detected
## by NGS at Site 1,393

- Agreement across platforms (high confidence)
  ≈ 17.6% (± 2.6%) minor component "A"

| Experiment | Reference "G" | Variant "A" | Coverage |
|---|---|---|---|
| EdgeBio PGM | 77.3% | 22.7% | 97 x |
| NIST PGM Run 1 | 82.5% | 17.5% | 2940 x |
| NIST PGM Run 2 | 83.4% | 16.6% | 3275 x |
| Illumina MiSeq | 83.7% | 16.3% | 26,234 x |
| Illumina HiSeq | 84.4% | 15.6% | 62,186 x |
| NIST SOLiD | 82.5% | 16.9% | 24,226 x |

Site 1,393 also confirmed by Niels Morling's lab using 454 technology (Martin Mikkelsen)

## Thanks for your attention!

Thanks to Tony Onorato and SWGDAM for the invitation to speak today

Questions and discussion?

Peter.Vallone@nist.gov
301-975-4872